

ICC: An Interconnect Controller for the Tofu Interconnect Architecture

August 24, 2010

Takashi Toyoshima

Next Generation Technical Computing Unit
Fujitsu Limited

shaping tomorrow with you

■ Requirements for Supercomputing Systems

■ Low latency

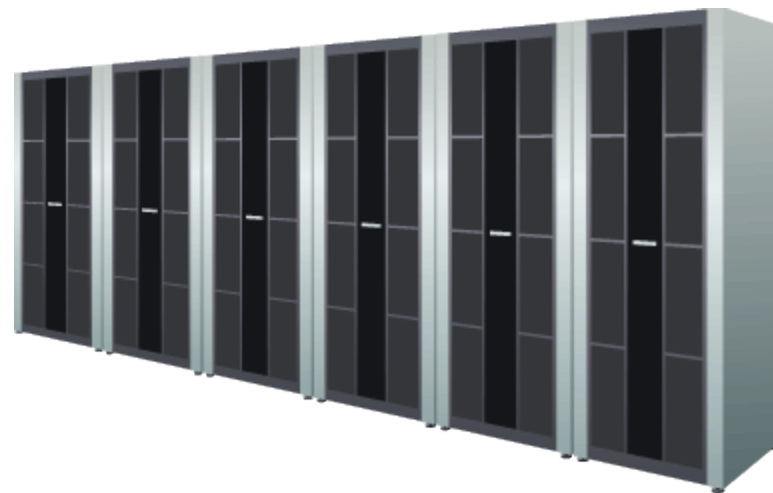
- Communication latency limits the scalability of applications

■ High bandwidth

- Increasing calculation FLOPS requires higher network bandwidth be balanced with FLOPS

■ RAS – Reliability, Availability and Serviceability

- The risk of hardware faults in large systems increases along with the increased number of nodes



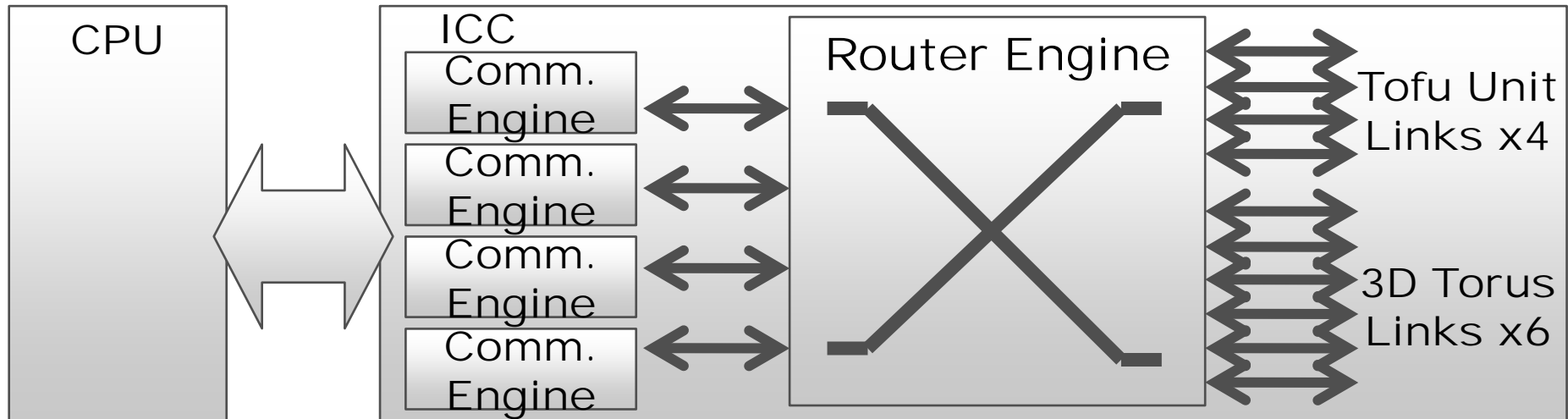
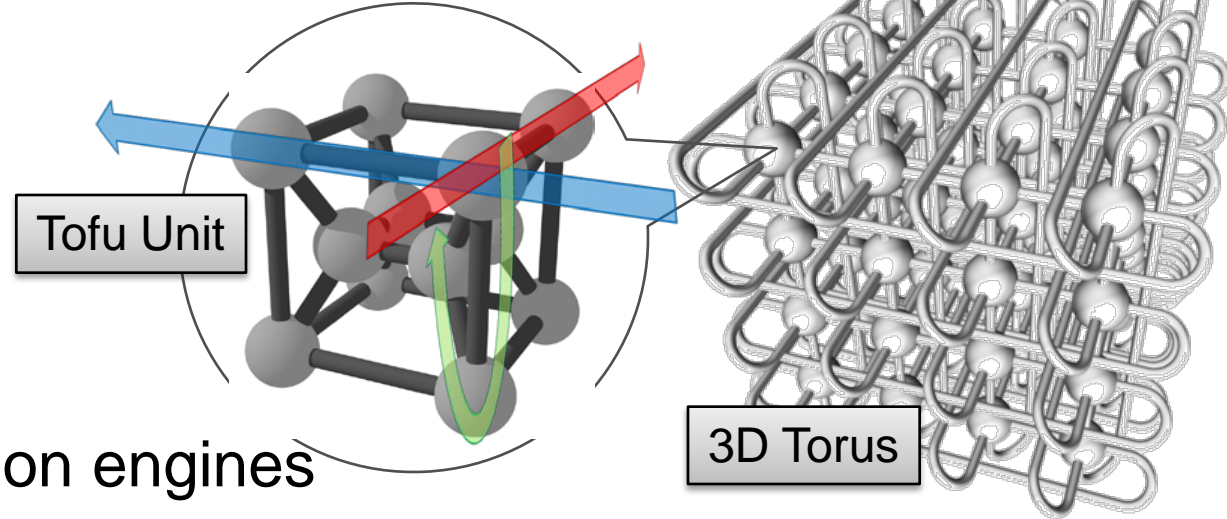
Fujitsu's New Interconnect Architecture

■ 6D Mesh/Torus Interconnect Architecture*

- Scalability
- Fault-tolerance

■ LSI Features

- Ten network links
- Four communication engines



(*) "Tofu: A 6D Mesh/Torus Interconnect for Exascale Computers", IEEE Computer, vol.42, no.11, Yuichiro Ajima, Shinji Sumimoto, Toshiyuki Shimizu

Implementation

- Implementation
- Features
 - Overview
 - Interface features for latency and throughput
 - Network features for network utilization
- Conclusion

■ Fujitsu's 65nm CMOS Technology

■ Die size

- 18.2mm × 18.1mm

■ Transistors

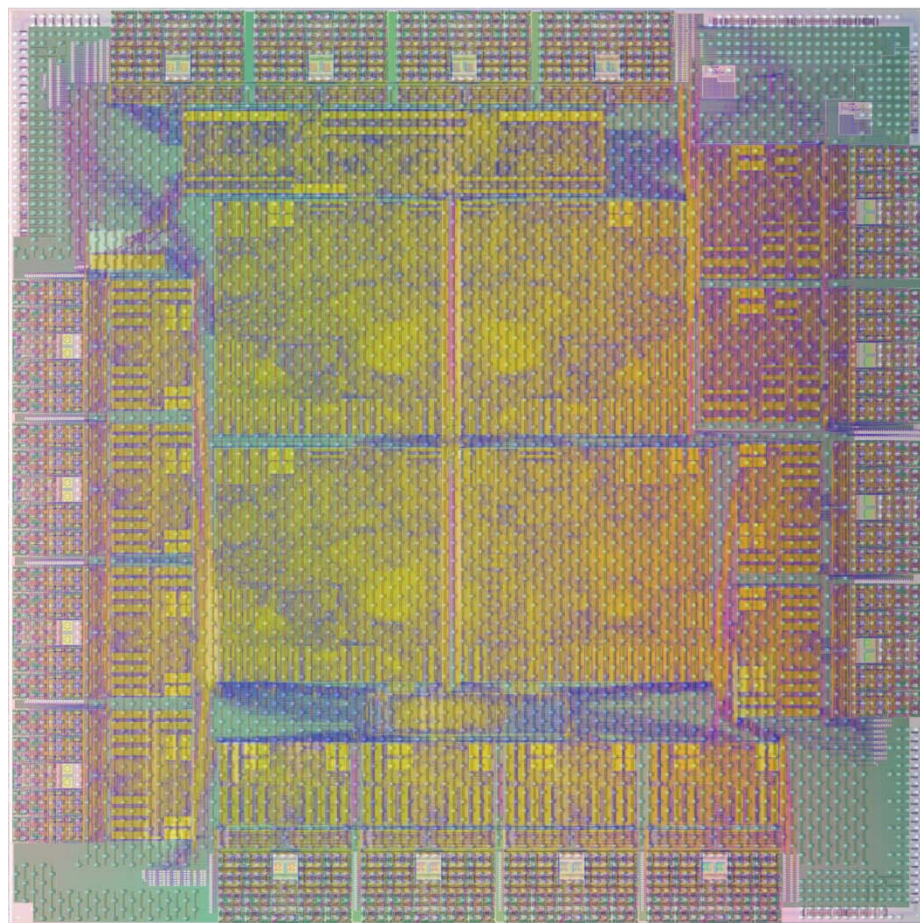
- 48M gates for logic
- 12M-bit SRAM cells

■ I/O

- 5GB/s Ports × 16
 - 6.25Gb/s × 8 links / port

■ Misc.

- ASIC design flow
- 312.5MHz/625.0MHz

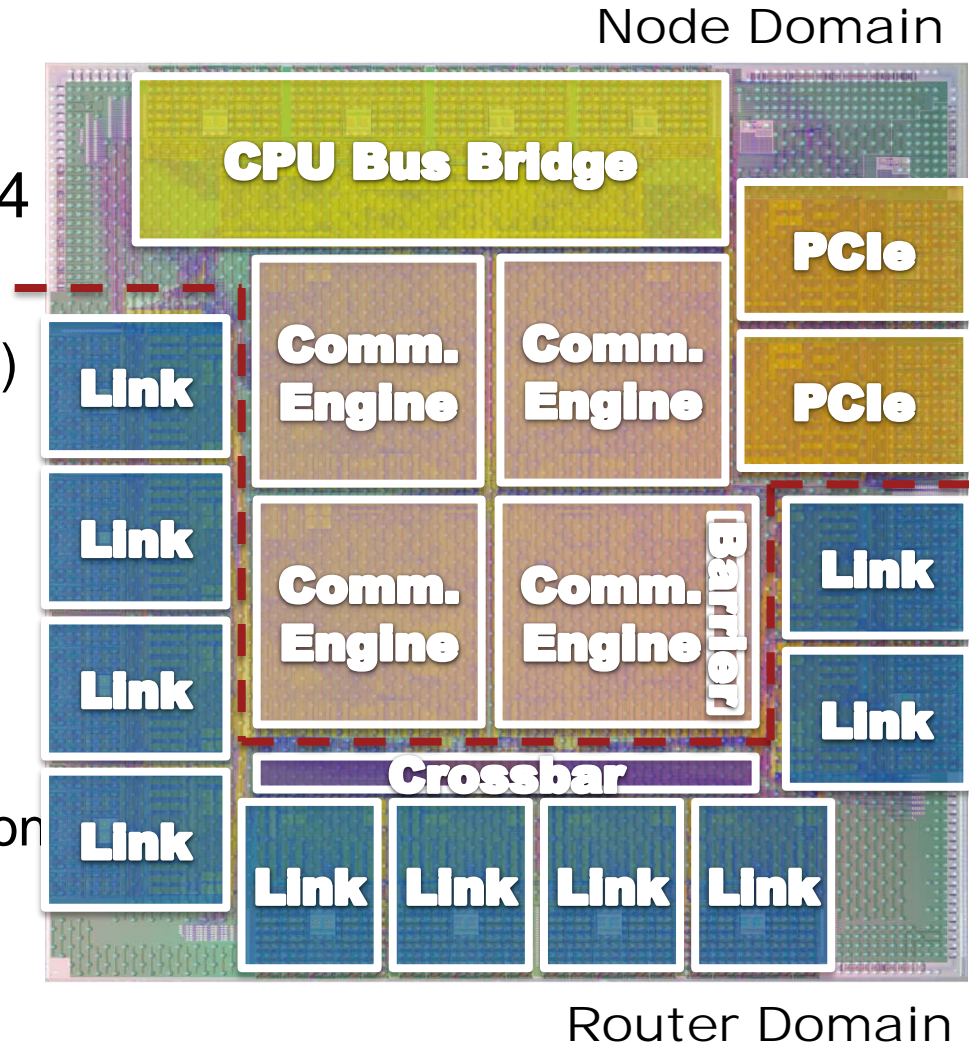


■ Node Domain

- CPU bus bridge
 - 20GB/s in each direction
- Communication engines × 4
 - 5GB/s in each direction
 - Barrier engine (Comm.#0 only)
- PCIe 2.0 root complex × 2
 - Isolated power domain

■ Router Domain

- Crossbar
 - 14 ports 5GB/s in each direction
- Link ports × 10
 - 5GB/s in each direction

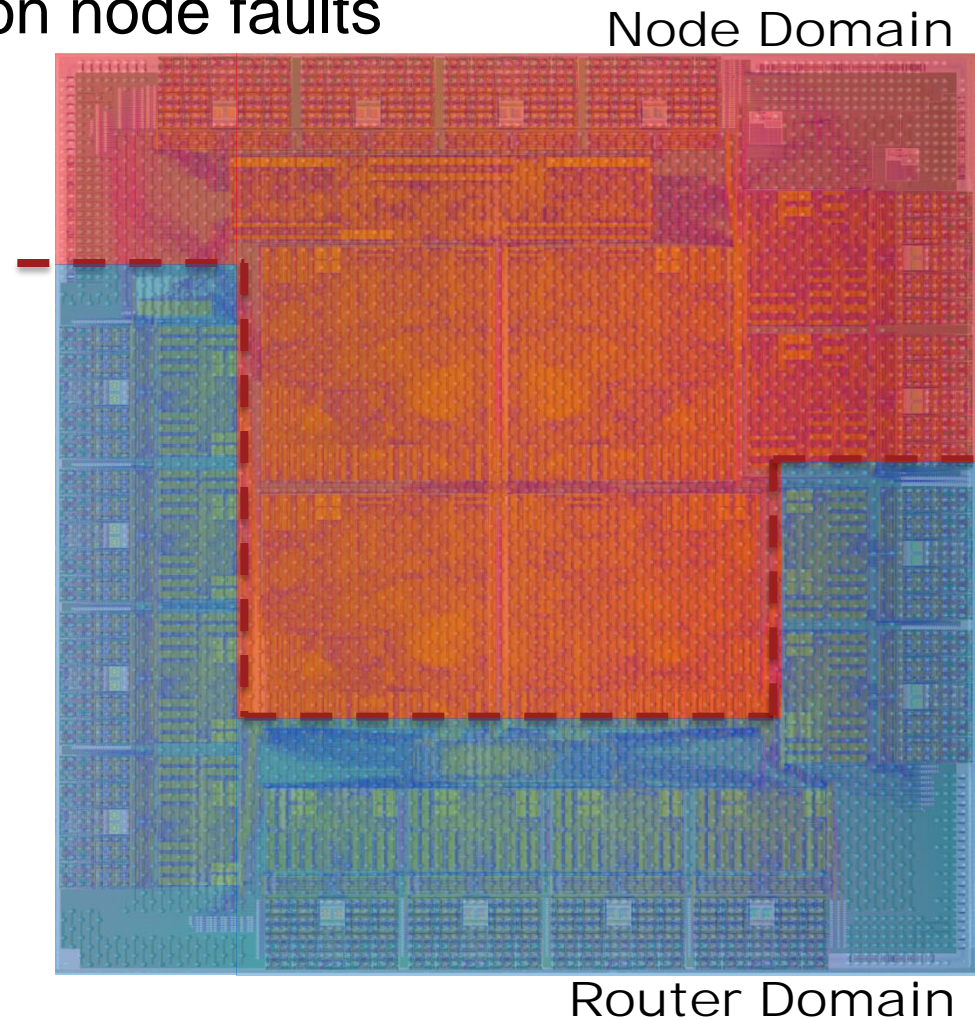


■ Fault Domain Isolation

- Router continues to work on node faults

■ Error Protection

- Radiation-hardened FFs
- ECC protection
 - RAM/Data path
- Parity error detection
 - Control path
- CRC protection
 - Data link/Transaction



Features Overview

- Implementation
- **Features**
 - **Overview**
 - Interface features for latency and throughput
 - Network features for network utilization
- Conclusion

	Latency	Throughput	RAS
System	<ul style="list-style-type: none"> ✓ Many Neighbors ✓ Hop Reduction ✓ 3D Torus View 	<ul style="list-style-type: none"> ✓ Many Neighbors ✓ Trunking 	<ul style="list-style-type: none"> ✓ Detour Path ✓ Subnet Partitioning
Network Interface	<ul style="list-style-type: none"> ✓ RDMA <ul style="list-style-type: none"> - Quick Start - Piggyback - Strong Order ✓ Stream Offload ✓ Barrier Engine 	<ul style="list-style-type: none"> ✓ GAP Control ✓ Multi-Interfaces <ul style="list-style-type: none"> - User Thread x2 - Kernel Thread 	<ul style="list-style-type: none"> ✓ Radiation-hardened FF ✓ ECC ✓ Parity ✓ CRC
Router Engine	<ul style="list-style-type: none"> ✓ Cut-through ✓ Grant Prediction ✓ Straight Bypass 	<ul style="list-style-type: none"> ✓ Straight Bypass ✓ New VC Scheduling 	<ul style="list-style-type: none"> ✓ Node Error Isolation ✓ Radiation-hardened FF ✓ ECC ✓ Parity ✓ CRC

★ : Unique Features

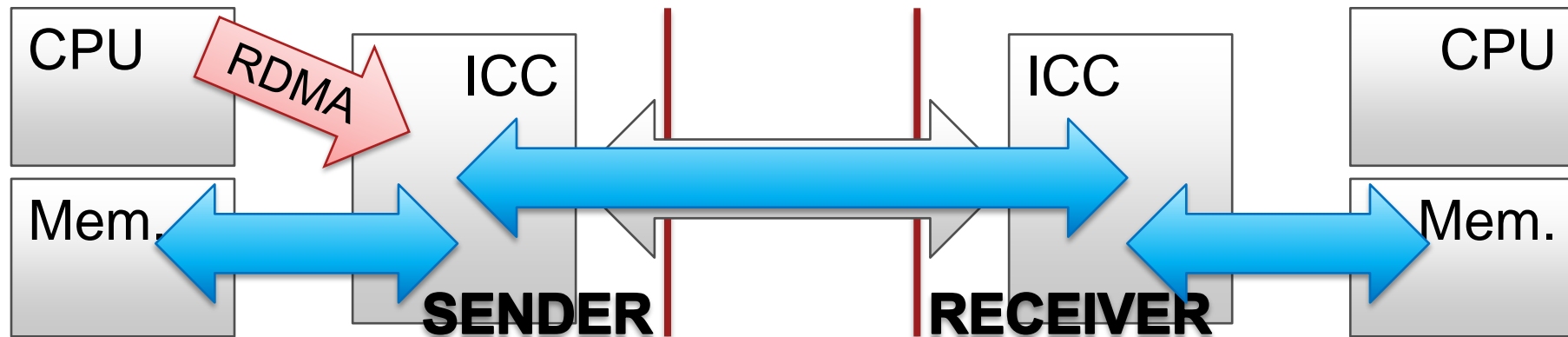
Today's topics are highlighted in red

Features

Interface features

- Implementation
- **Features**
 - Overview
 - **Interface features for latency and throughput**
 - Network features for network utilization
- Conclusion

■ Features



Operation	READ (Get) / WRITE (Put)
Length	~16MB
MTU	256B~1920B
Virtual Address	Support (64K set)

■ Low Latency and High Throughput

- Command supply throughput and latency
- Out-of-ordered I/O memory bus

■ Sender Techniques

■ Direct descriptor

- Quick command supply

	Throughput	Latency
PIO		✓ Good
DMA	✓ Good	

Command Supply Performance

■ Piggyback

- Command embedded communication payload
- Short message sending without any DMA

■ Receiver Techniques

■ Out-of-ordered I/O memory bus

- High throughput bus transaction

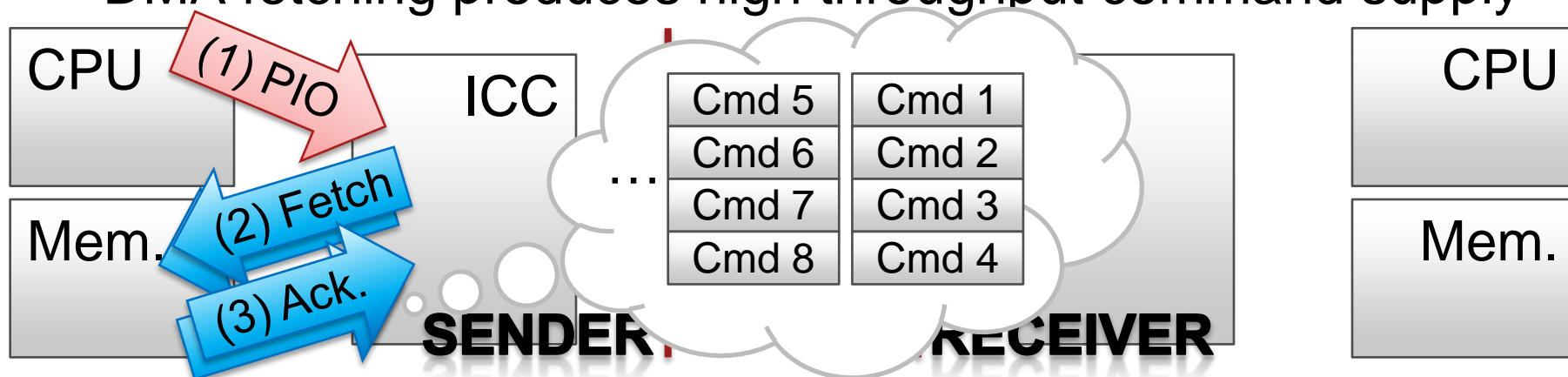
■ Strong ordered store

- In order completion of DMA transactions for buffer polling

Direct Descriptor Feature

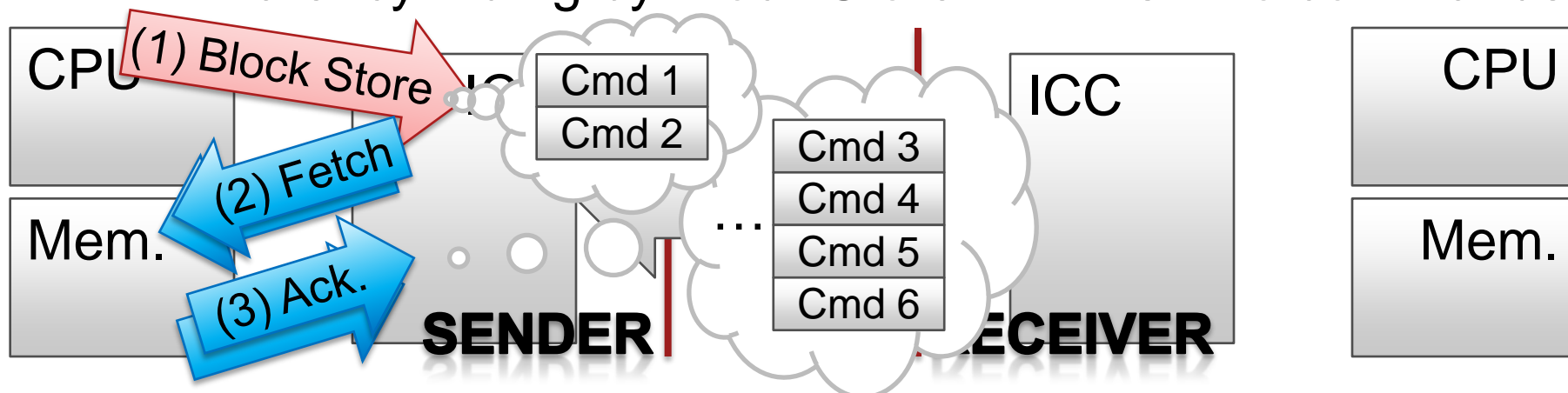
■ Normal Command Supply

- DMA fetching produces high throughput command supply



■ Direct Descriptor and DMA Command Supply

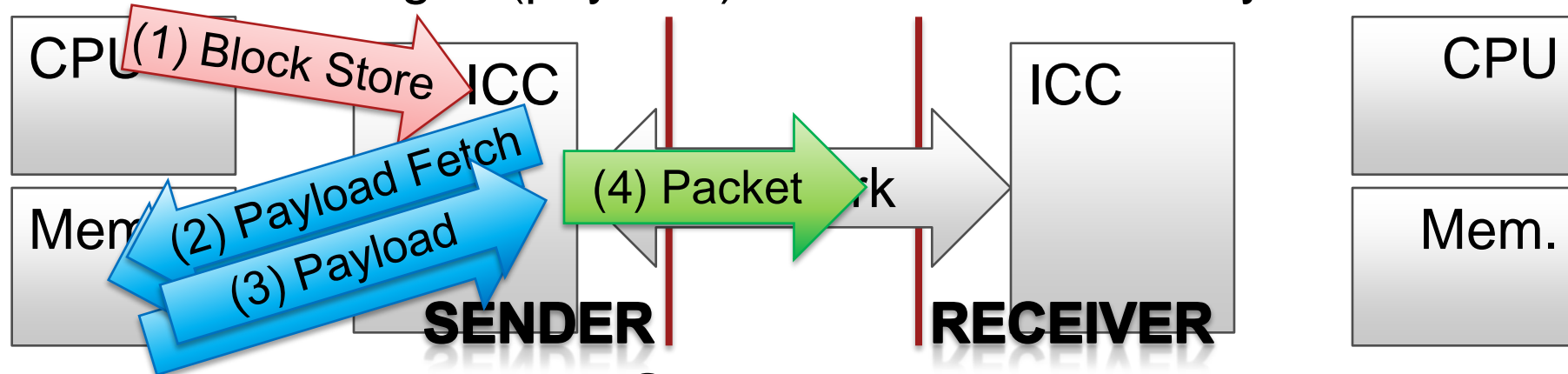
- DMA latency hiding by Block Store with first two commands



Piggyback Feature

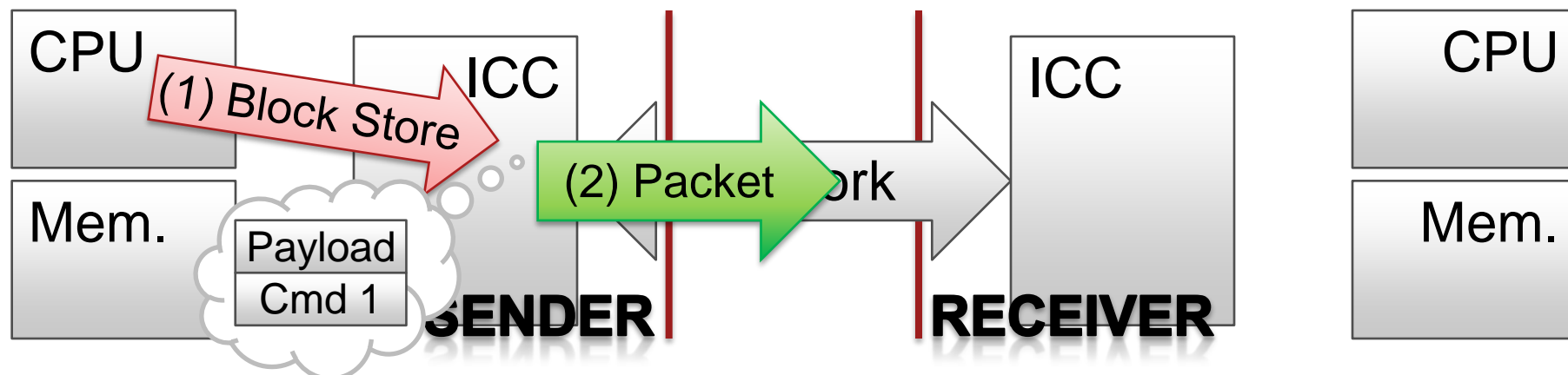
■ Normal Payload Supply

- User messages (payload) should be fetched by DMA



■ Piggyback Payload Supply

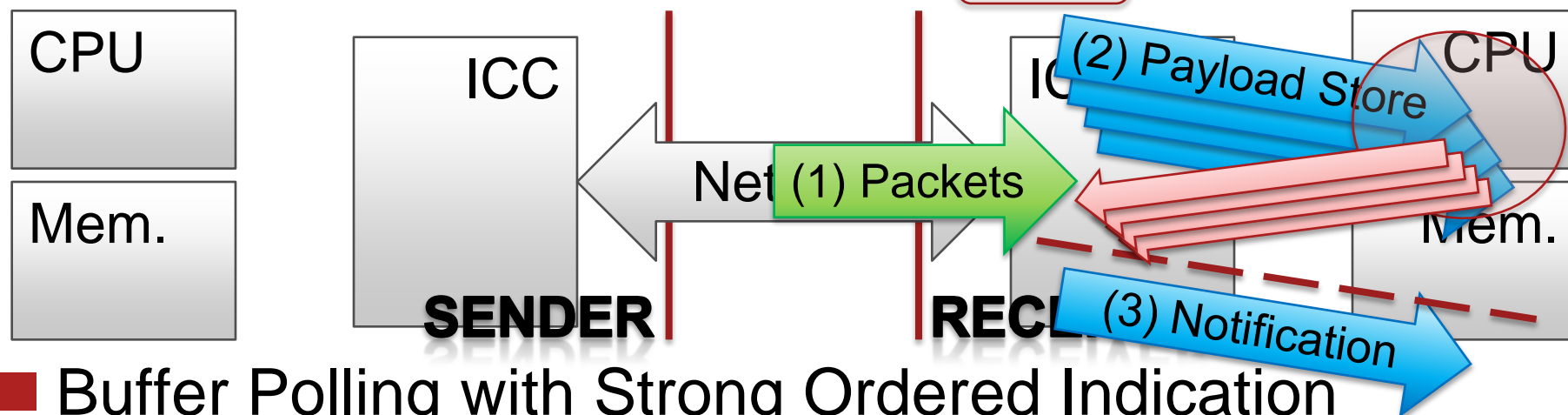
- User messages (payload) are embedded in commands



Out-of-Ordered I/O Memory Bus

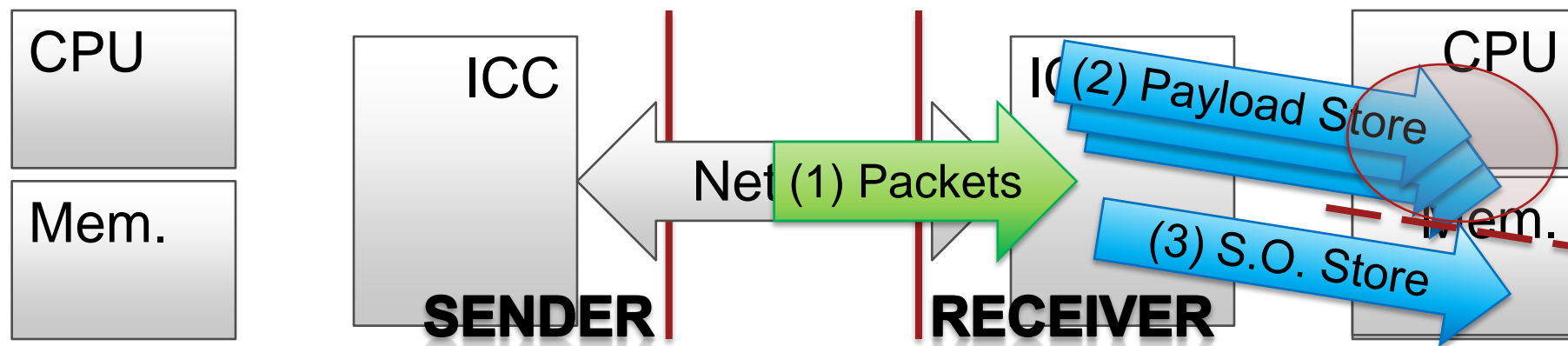
■ Completion Notification Polling

- ICC notifies after the completion of **O-o-O** DMA Stores



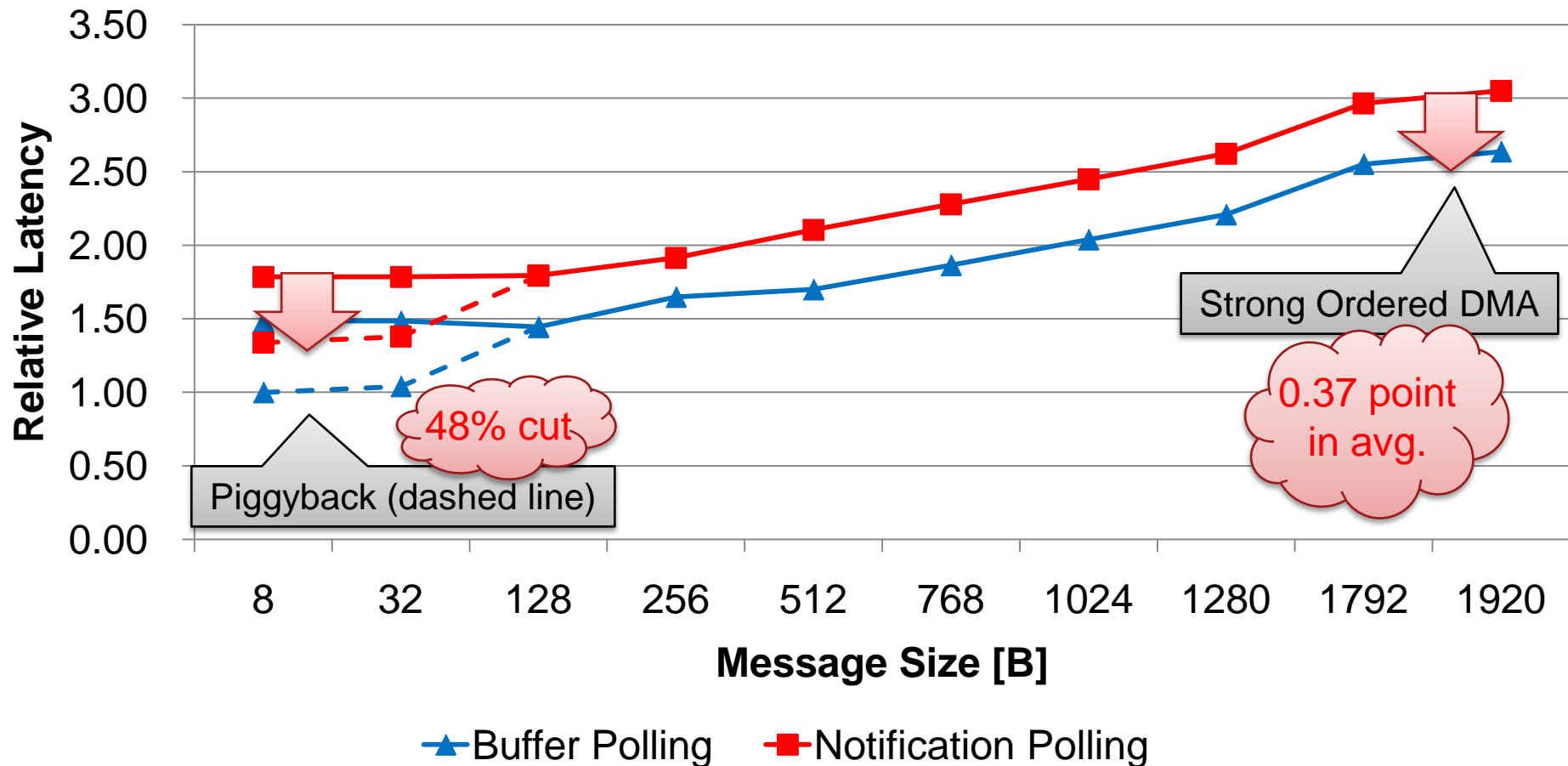
■ Buffer Polling with Strong Ordered Indication

- Memory controller guarantees specified DMA ordering



■ Hardware Measured Results

- Piggyback achieves low latency in short message
- Strong ordered packet makes buffer polling possible



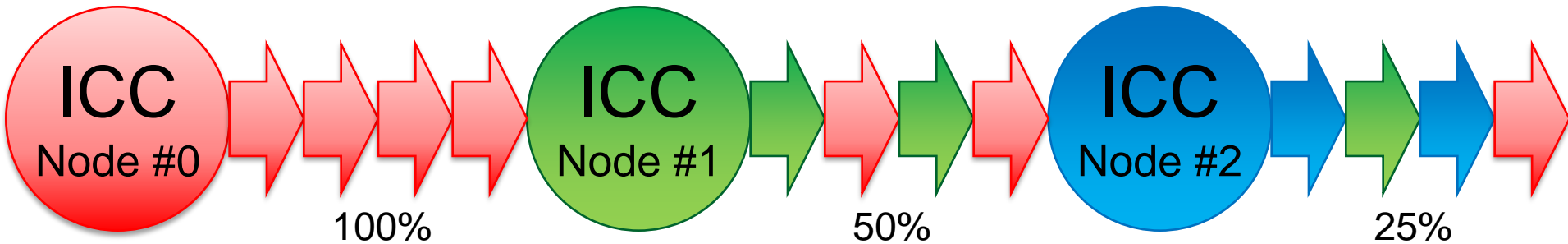
Features

Network Features

- Implementation
- **Features**
 - Overview
 - Interface features for latency and throughput
 - **Network features for network utilization**
- Conclusion

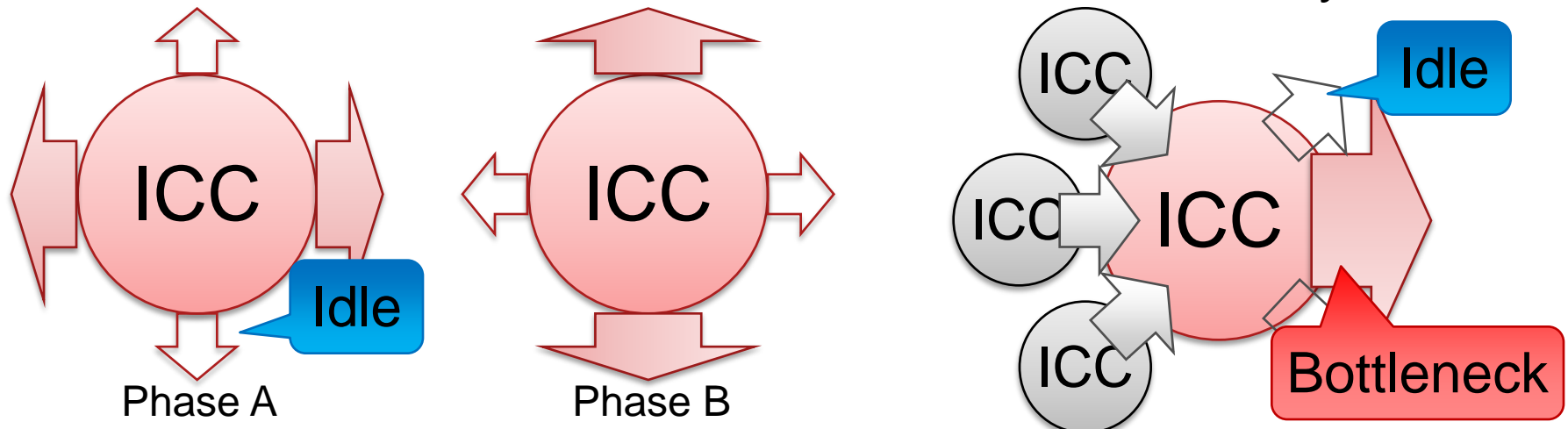
■ Global Unfairness of Throughput

- Arbitrations with local fairness cause global unfairness



■ Non-uniform Application Traffic in Time and Space

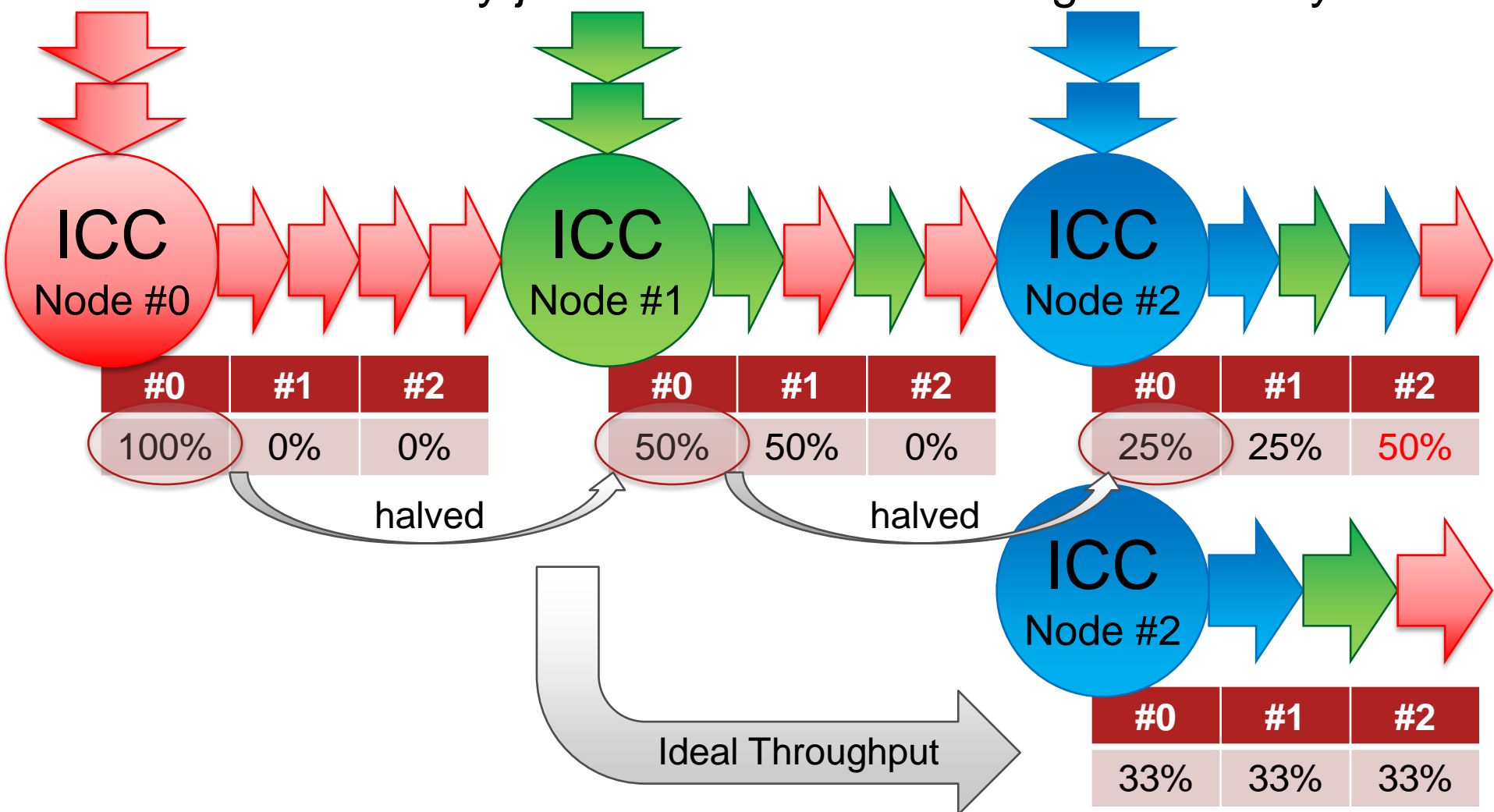
- Bandwidth of idle links needs to be used effectively



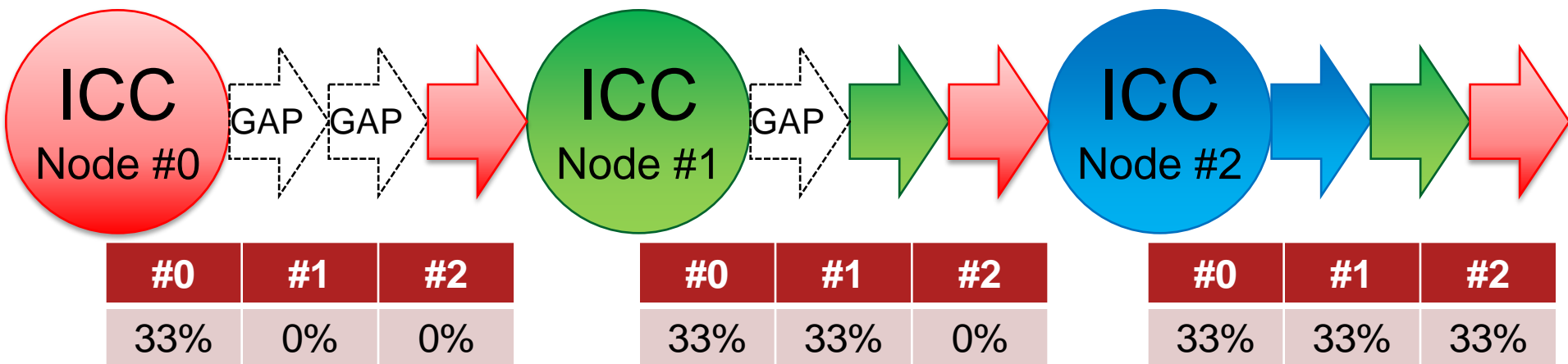
Global Unfairness of Throughput

Local Fairness of arbitration cause global unfairness

- Arbiters on every junction treat all incoming traffic fairly

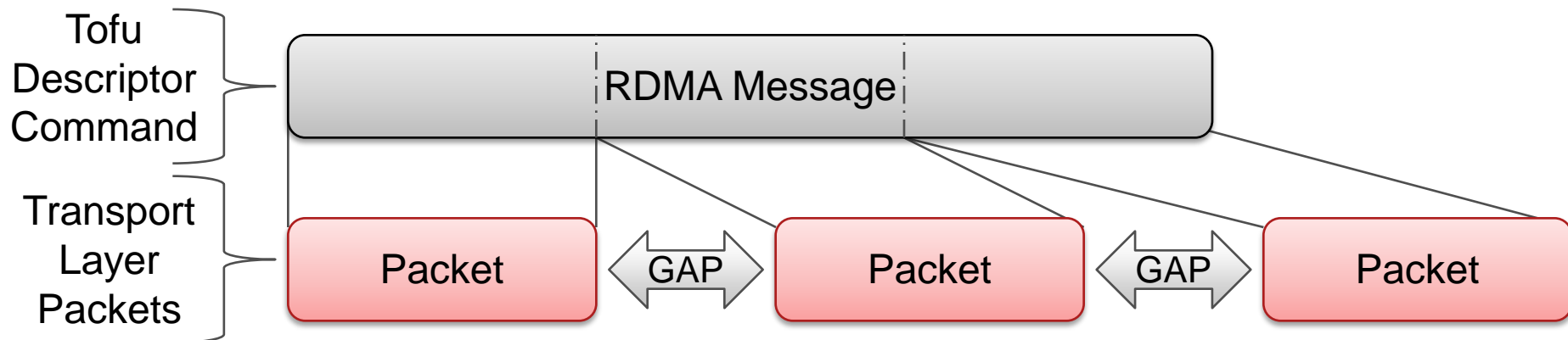


■ Software Specify the Inter-Packet GAP Parameter



■ Communication engine works to control injection rate

- Insert temporal gaps between transmitting packets
- Interval can be specified by the user

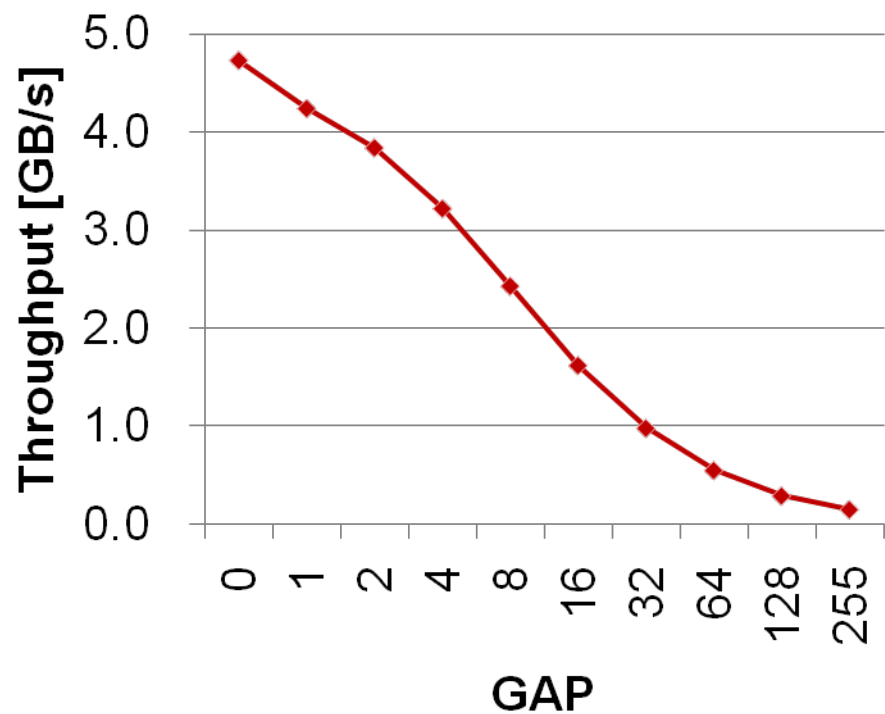


Throughput Performance

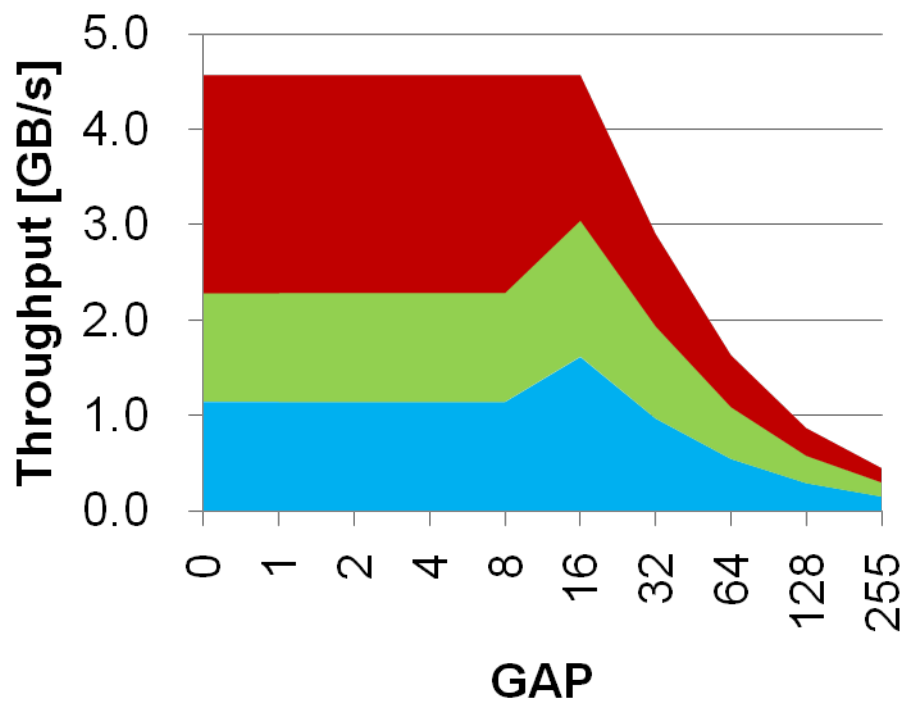
Hardware Measured Results

- Software can specify fine grained GAP parameters: 0-255
- GAP works to control throughput effectively

GAP Sensitivity



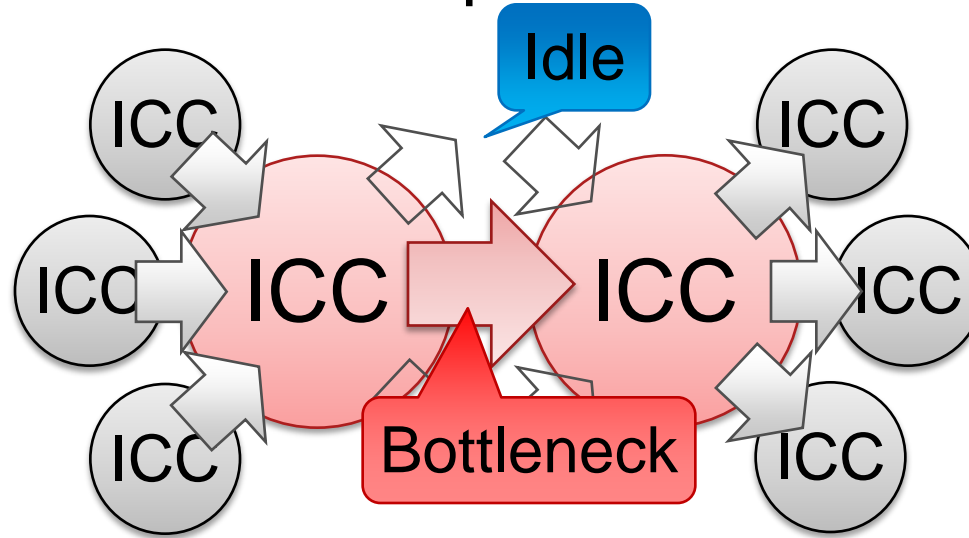
Stacked Throughput



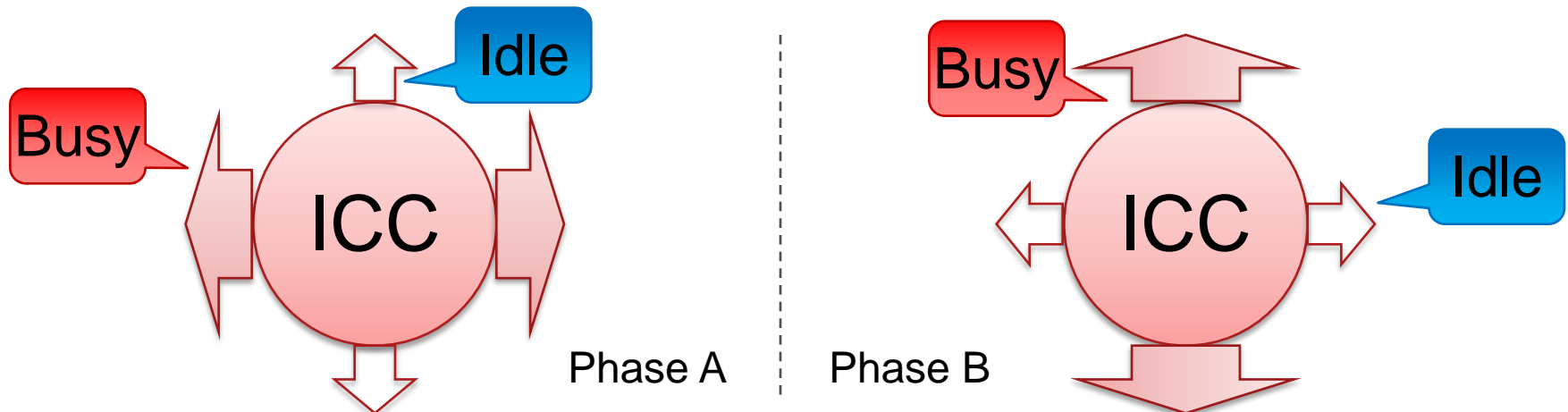
Node #0 Node #1 Node #2

Non-uniform Application Traffic

■ Non-uniform Traffic in Space



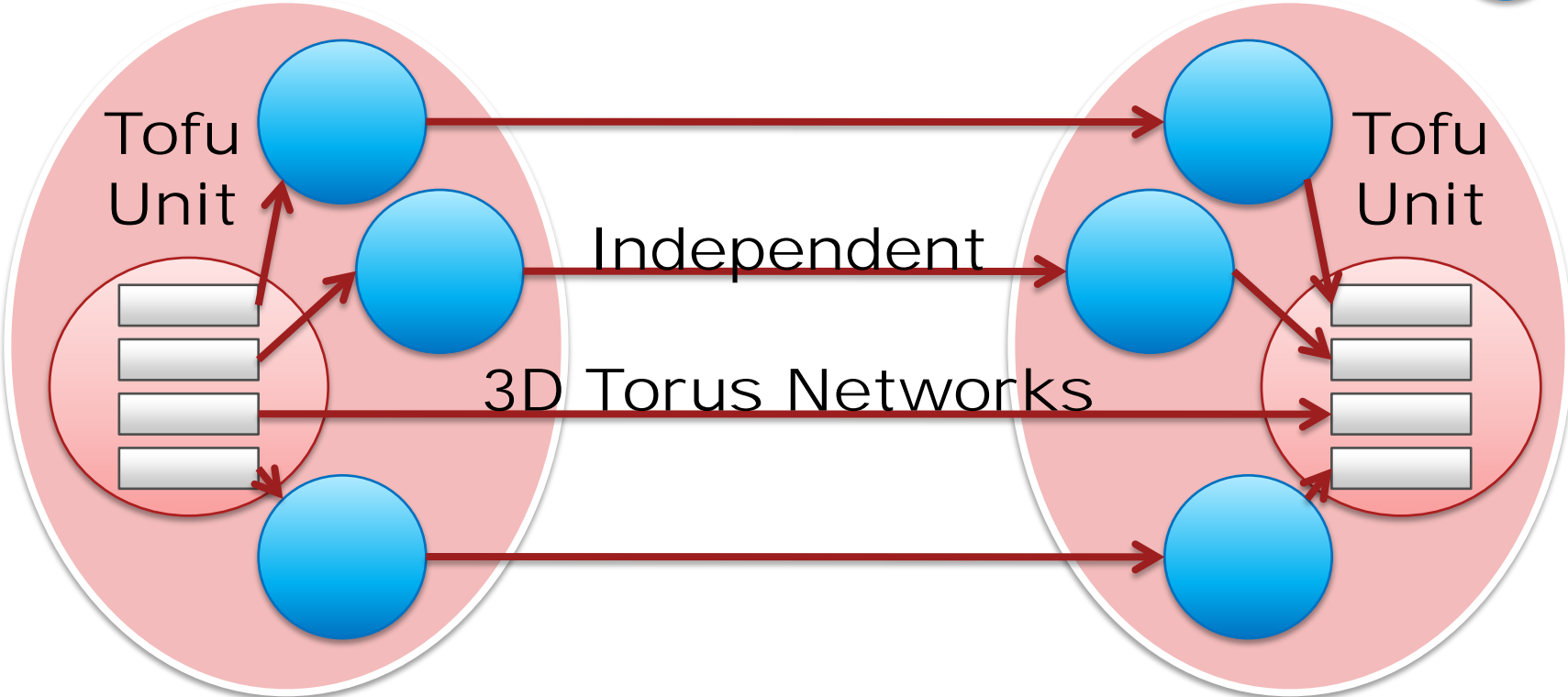
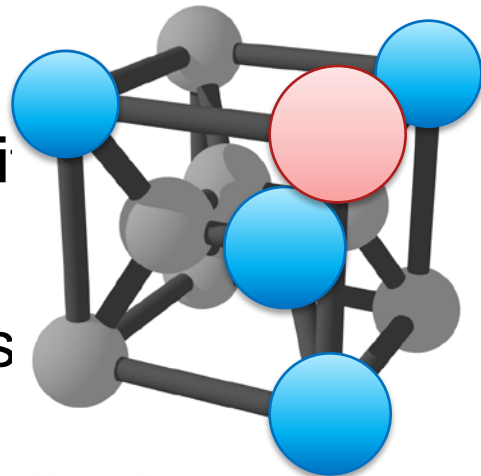
■ Non-uniform Traffic in Time



Trunking Communication

■ Trunking Independent Idle Paths

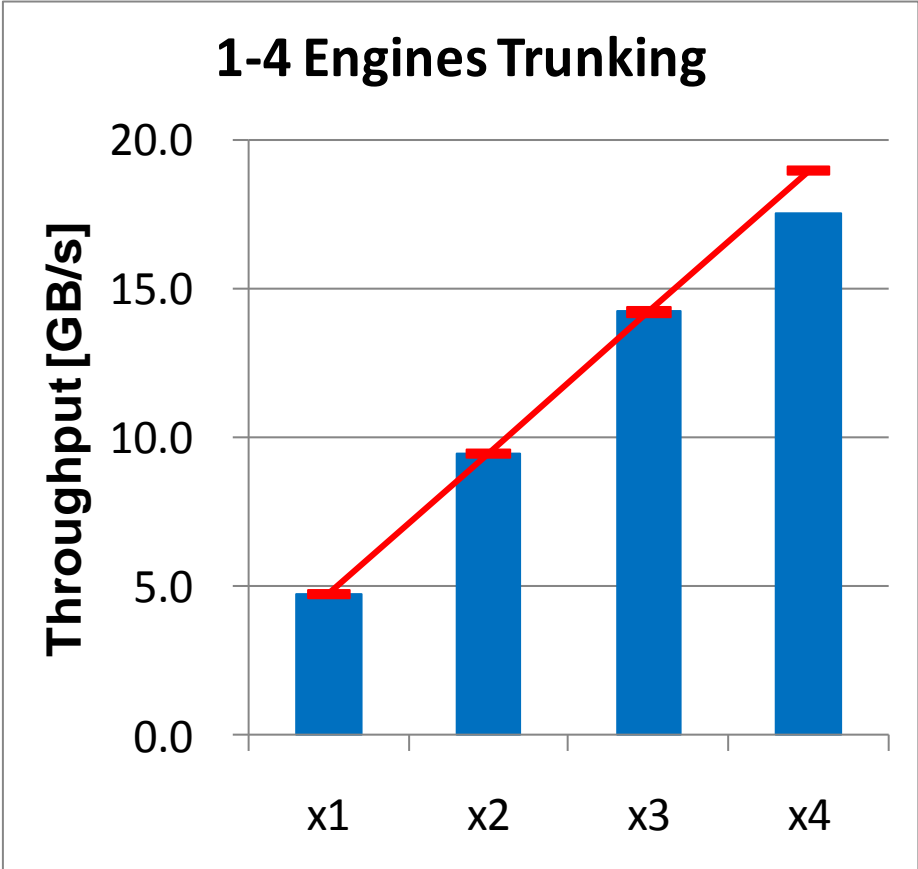
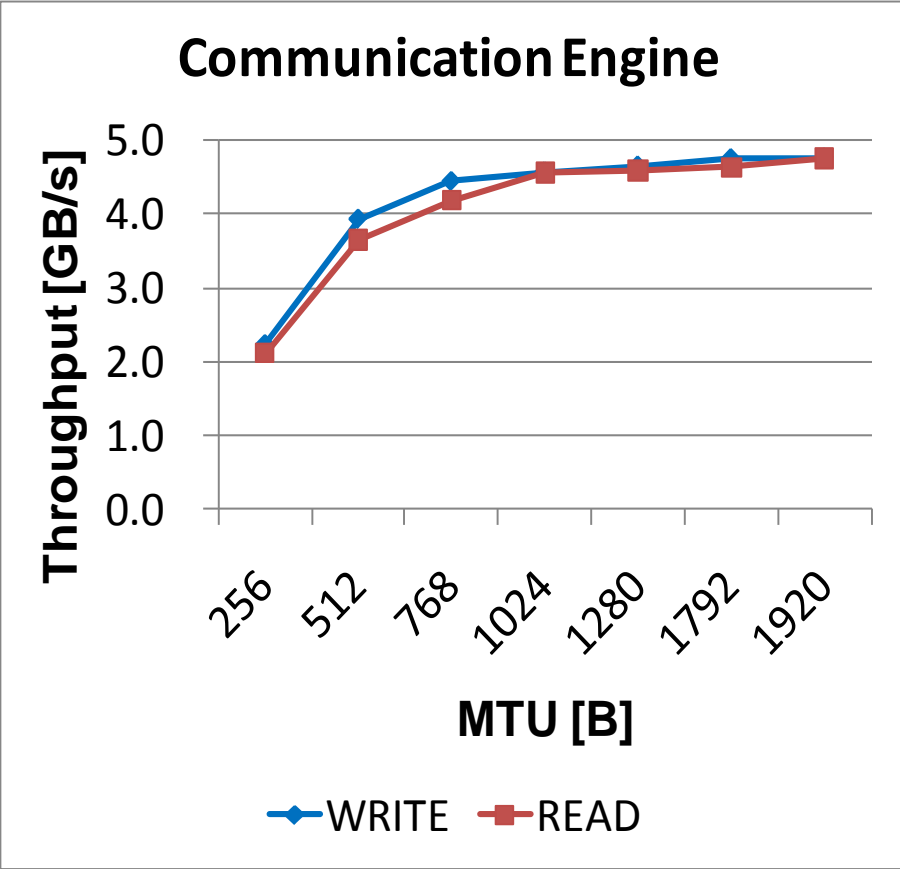
- Nodes have four neighborhoods in Tofu Uni
 - Independent links and 3D-Torus networks
- Each node has four communication engines
 - Up to $\times 4$ throughput



Trunking Performance

Results

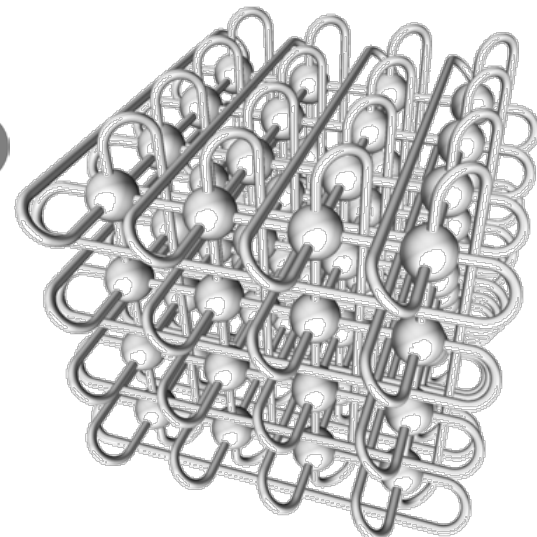
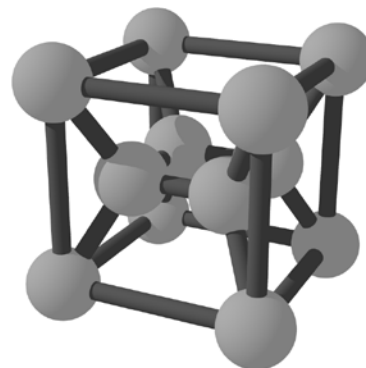
- Communication engines achieve good performance
- Trunking mechanisms scale up to four engines



Conclusion

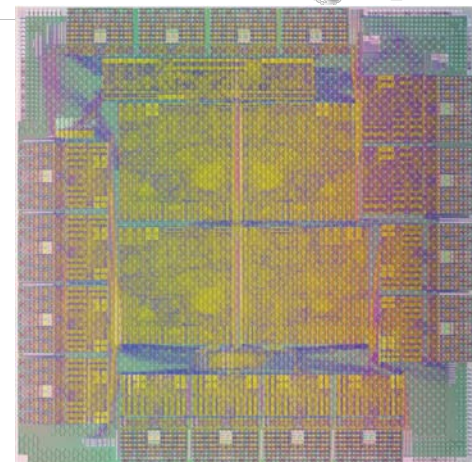
- Implementation
- Features
 - Overview
 - Interface features for latency and throughput
 - Network features for network utilization
- Conclusion

- Tofu: A 6D mesh/torus interconnect architecture
 - Interconnect for Fujitsu's Peta/Exascale computing systems
 - Low latency, High bandwidth and RAS



■ Features

- High-throughput and low-latency RDMA
 - Direct Descriptor and Piggyback
 - Out of Order I/O Memory Bus
- Network features for network utilization
 - Network injection rate control
 - Trunking up to four times throughput



Thanks to...

Next Generation Technical Computing Unit

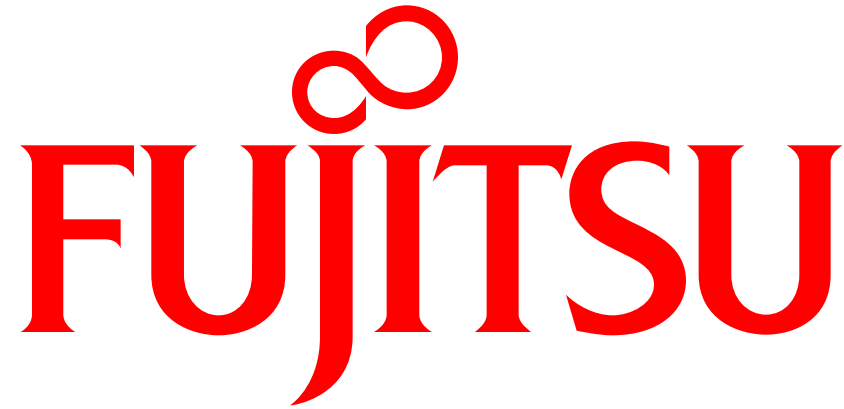
Aiichiro Inoue, Yuji Oinaga

Tofu Architecture Team

Toshiyuki Shimizu, Yuichiro Ajima, Tomohiro Inoue, Shinya Hiramoto

ICC Design Team

Takeo Asakawa, Akira Asato, Takumi Maruyama,
Koichiro Takayama, Koichi Yoshimi, Osamu Moriyama,
Masao Yoshikawa, Shinichi Iwasaki, Takekazu Tabata,
Yoshiro Ikeda, Yuzo Takagi,
Yoshihito Matsushita, Toshihiko Kodama, Satoshi Nakagawa,
Masato Inokai, Shigekatsu Sagi, Ikuto Hosokawa,
Yaroku Sugiyama, Takahide Yoshikawa



shaping tomorrow with you